

Lecture 9

2 layer network

sol 2 - NEW

sol 1 - PM

OH - 9:30, 10:30

[HW4 - TB part 2]

Lecture Notes IV – Neural Networks, Part 1

Marina Meilă
mmp@uwaterloo.ca

With Thanks to Pascal Poupart & Gautam Kamath
Cheriton School of Computer Science
University of Waterloo

February 2, 2026

A little history ✓

The single "neuron" ✓

Two-layer Neural Networks ←
Hidden layer options
Output layer options

Multi-layer neural networks ←

Training a neural network by backpropagation ↪...

Reading HTF Ch.: 11.3 Neural networks, Murphy Ch.: (16.5 neural nets), Bach Ch.: –, Deep Learning Book (Goodfellow, Bengio, Courville) 6.1-4, ResNet 7.6, ConvNet 9., Autoencoders 14.1, Dive Into Deep Learning 4.1-4.3.

(Artificial) Neural Network (nn) unit

- ▶ For each **unit** i

$$y_i \equiv f_i(x) = \phi\left(\sum_j w_{ij}x_j + w_{i0}\right) \quad (1)$$

- ▶ **Weigth vector** w_i

- ▶ w_{ij} = strength of the link from unit i to input j
- ▶ $w_{ij} = 0$: no link
- ▶ w_{ij} can be positive or negative
- ▶ Sometimes we call the input vector $x = [x_{1:d}]$ **input units**

- ▶ **activation function** $\phi()$

- ▶ must be non-linear (otherwise the unit is a linear transformation)
- ▶ wanted: monotonically increasing, differentiable, gradient non-zero¹

Notation ϕ is overloaded

non-linear, increasing

- ▶ When we talk about nn in general: ϕ is any activation function
- ▶ When we do calculations with nn: ϕ is by default the **logistic** function (unless specified otherwise)

logistic or sigmoid function

$$\phi(u) = \frac{1}{1 + e^{-u}}$$

(2)

- ▶ When we do statistics or ML (but not nn): ϕ is the logistic function
- ▶ Exercise: compare $f_i(x)$ from (1) with $p(x) = Pr[y = 1 | x]$ from logistic regression.

¹More technically: ϕ can be any continuous, bounded and strictly increasing function on \mathbb{R} .

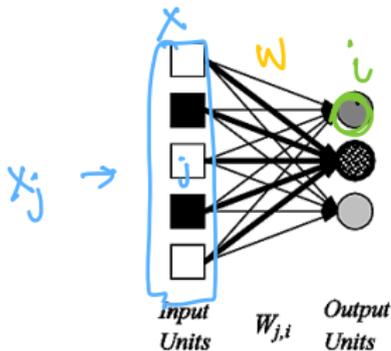
$$\mathbf{x} \rightarrow \underline{\underline{\beta^T \mathbf{x}}} = z = f(x)$$

$$P(y|x, \beta) = \varphi(yz)$$

$$= \frac{1}{1 + e^{-yz}}$$

Perceptron

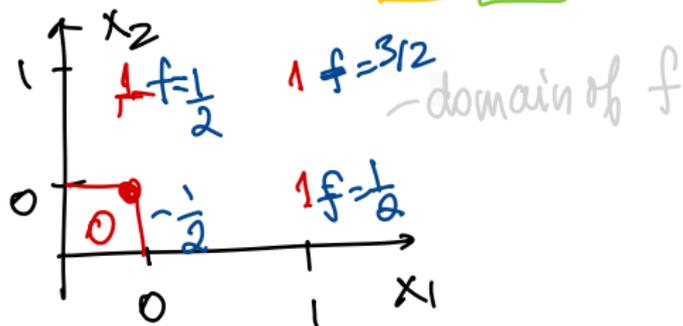
- Single layer feed-forward network



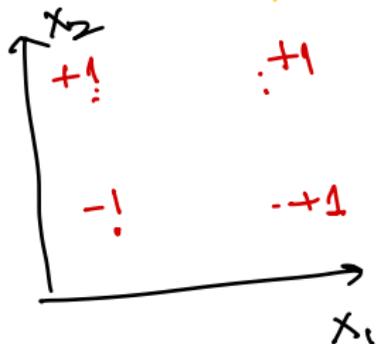
$$\varphi(z) = \frac{1}{1 + e^{-z}}$$

logistic sigmoid

Perceptron weights for OR, AND, XOR \rightarrow Exercise



OR as classification



$$f^{true}(x_1, x_2) = x_1 \text{ OR } x_2$$

$$x_{1,2} \in \{0, 1\}$$

$$f(x_1, x_2) = w_1 x_1 + w_2 x_2 + w_0$$

$$w_0 = -1/2$$

$$w_1 = 1$$

$$w_2 = 1$$

$$y = \frac{\text{sgn}(f(x_1, x_2)) + 1}{2}$$

try logistic regression

$$w_{1,2} = 1$$

$$w_0 = -1/2$$

(linear output)

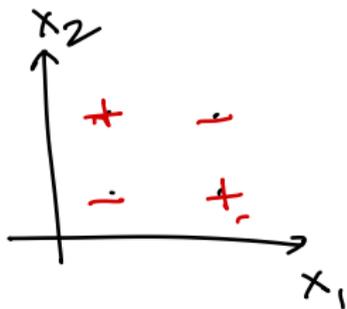
ϕ = activation

confidence
 $P(y|x) =$

$$\phi(y|x)$$

\uparrow
 logistic



Perceptron weights for OR, AND, XOR

$$f^{\text{true}}(x_1, x_2) = \frac{x_1 \text{ XOR } x_2 + 1}{2}$$

$w_0, w_1, w_2 = ? \quad \exists w_{0,1,2}$ so that $\text{sgn}(w_0 + w_1 x_1 + w_2 x_2) = y_i$ for all i

get around

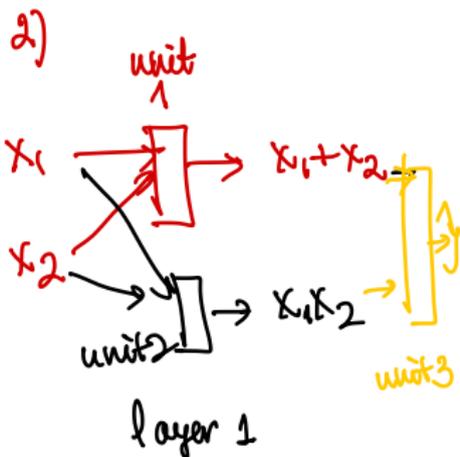
$$1) \quad \tilde{x} = [x_1 \quad x_2 \quad x_1 x_2 \quad 1]$$

$$\tilde{w} = [w_1 \quad w_2 \quad w_3 \quad w_0]$$

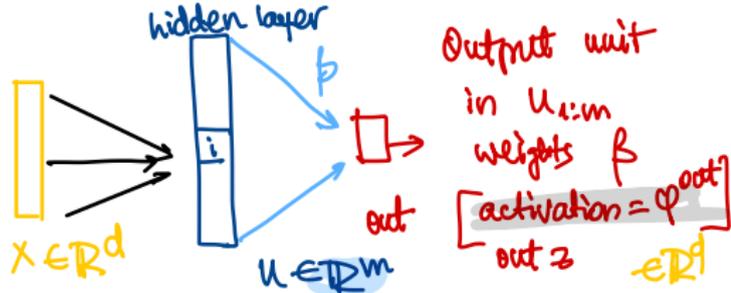
intercept
as slope

$$w_1 = w_2 = 1 \quad w_0 = -\frac{3}{4}$$

$$w_3 = -\frac{3}{2}$$



Two-layer Neural Networks



- We build a **two-layer neural network** in the following way:

Inputs	x_j	$j = 1 : d$
Bottom layer ²	$u_i = \phi(w_i^T x)$	$i = 1 : m, w_i \in \mathbb{R}^d$
Top layer	$f = \beta^T u$	$\beta \in \mathbb{R}^m$
Output	f	$\in \mathbb{R}$

unit i : input x
weights w_i
activation ϕ
out u_i

In other words, the neural network implements the function

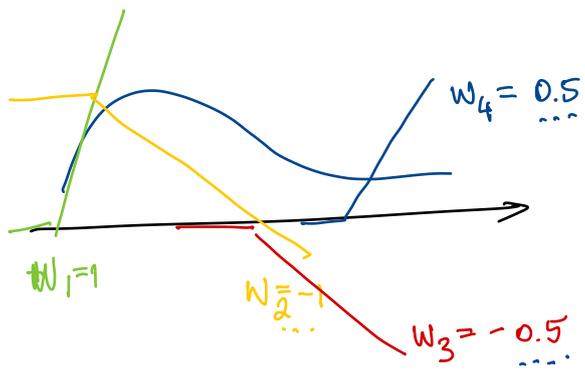
$$f(x) = \sum_{i=1}^m \beta_i u_i = \sum_{i=1}^m \beta_i \phi\left(\sum_{j=1}^d w_{ij} x_j\right) \in (-\infty, \infty) \quad (3)$$

Annotations: z (red), $linear$ (blue), $linear$ (yellow), $z = \beta^T u$ (red)

Note that this is just a linear combination of logistic functions.

- As we will see shortly, in general, $f(x)$ can also be non-linear

²In neural net terminology, each variable u_i is a **unit**, the bottom layer is **hidden**, while top one is **visible**, and the units in this layer are called hidden/visible units as well. Sometimes the inputs are called **input units**; imagine neurons or individual circuits in place of each x, u, y variable.



Output layer



- ▶ Let $z = \beta^T u$ be the **linear output** of the nn. **hidden**
- ▶ For problems other than regression, it unifies the analysis to apply a non-linear function ϕ_{out} to z . This is why, **theoretically** but not in practice, we will write $f(z(x)) = \phi_{\text{out}}(z(x))$
- ▶ Why? $\phi_{\text{out}}(z)$ matches the loss \mathcal{L} , which in turn matches the **prediction problem**

Prediction problem	Predict	Output layer ϕ_{out}	Range of f
→ regression	$\hat{y} = z$	linear $\phi_{\text{out}} = z$	$\in \mathbb{R}$
→ <u>binary classification</u>	$\hat{y} = \text{sgn} z$	logistic $\phi_{\text{out}} = \phi(\beta^T u)$	$\in [0, 1] \leftarrow \text{logistic}$
→ <u>multiway classification</u>	$\hat{y} = \underset{1:r}{\text{argmax}} z_k$	softmax $\phi_{\text{out}k} = \phi_k(\beta_k^T u)$	$\in [0, 1]^r$

$$z_k = \beta_k^T u \quad k=1:r$$

- ▶ Regression: **linear** layer as in (3) $f = \sum_i \beta_i u_i$
- ▶ Classification (binary): **logistic** layer $f(x) \in [0, 1]$ is interpreted as the probability of the + class.

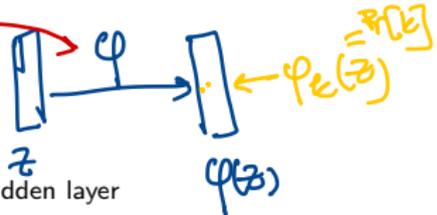
$$f(x) = \phi \left(\sum_{j=1}^m \beta_j u_j \right) = \phi \left(\sum_{i=1}^m \beta_j \phi \left(\sum_j w_{ij} x_j \right) \right) = \text{Pr}[y=+ | x] \quad (4)$$

- ▶ Multiway classification with r classes
 - ▶ Output is vector of r functions f_1, \dots, f_r
 - ▶ f_k is the probability of $y = k$
 - ▶ (sometimes f_k can be a "confidence")
 - ▶ This is done with a **softmax** layer (next page)

$$p_k = \text{Pr}[y=k | x]$$

The softmax function

- ▶ logistic $\phi(z) = \frac{e^z}{1+e^z}$
 - ▶ represents the probability that $y = 1$ in binary classification
 - ▶ in a nn, $f(x) = \phi(\beta^T u)$, with $u \in \mathbb{R}^m$ the activations of the hidden layer
- ▶ The **softmax** function generalizes the logistic to r classes
 $\phi(z) : \mathbb{R}^r \rightarrow (0, 1)^r$



$$[z_1, \dots, z_n]$$

$$\phi_k(z) = \frac{e^{z_k}}{\sum_{j=1}^m e^{z_j}}, \text{ for } k = 1 : r$$

$$z_k \rightarrow e^{z_k} \text{ normalize} \quad (5)$$

$$\phi(z) = [\phi_1(z) \dots \phi_r(z)] \text{ (overload of } \phi)$$

(6)

Properties of softmax

- ▶ $\sum_{j=1}^m \phi_j(z) = 1$ for all z
- ▶ for $z_k \gg z_j$, $\phi_k(z) \rightarrow 1$.
- ▶ derivatives $\frac{\partial \phi_j}{\partial z_k} = \phi_k \delta_{jk} - \phi_j \phi_k$
- ▶ in a nn, with output activations $u \in \mathbb{R}^m$, we train r weight vectors $\beta_{1:r}$, and

$$f_k(x) = \phi_k(\beta_k^T u(x)), \quad \text{for } k = 1 : r \quad (7)$$

Are multiple layers necessary?

- ▶ 1990's: NO
- ▶ 2000's: YES
- ▶ 2020's: The more the better!
- ▶ A theoretical result

Theorem (Cybenko, \approx 1986)

Any continuous function from $[0, 1]^d$ to \mathbb{R} can be approximated arbitrarily closely by a linear output, two layer neural network defined in (3) with a sufficiently large number of hidden units m .

- ▶ A practical result



Deep Learning

Deep learning = multi-layer neural net

- ▶ So, what is new?
 - ▶ small variations in the "units", e.g. switch stochastically w.p. $\phi(\mathbf{w}^T \mathbf{x}^{in})$ (Restricted Boltzmann Machine), Rectified Linear units
 - ▶ training method stochastic gradient, auto-encoders vs. back-propagation (we will return to this when we talk about training predictors)
 - ▶ lots of data
 - ▶ **double descent**

